Volker Roth

# Rechnersicherheit, SoSe 21
## Übung 03
**TutorIn: Oliver Wiese**
**Tutorium 02**
**Materialien: Latex, VSC, Skript**

19. Mai 2021

---

## 1 Password models

(a) *Choose your tool to analyze the password set. We recommend python (with pandas) or R but you can choose an other tool. Be aware that most CSV-Tools can not handle a huge file.*

- Wir haben das pandas CSV-Tool benutzt. Da die zu analysierende Datei jedoch keinen seperator zwischen Passworthäufigkeit und Passwort hat, haben wir folgendes versucht:

  1. At first we used a white space as separator, however, this caused some passwords to be left out or incorrectly recorded, which threw several errors as consequence.

     ```
     df = pd.read_csv(p_datafile, sep=' ', header=None,  names=['frequentcys', '
         passwords'])
     ```

  2. After that we used a regex as separator, however, the number of errors only decreased slightly. It also took a lot longer to analyze the file.

     ```
     sep = r'^(\s*[0-9]+)'
     df = pd.read_csv(p_datafile, sep, header=None,  names=['frequentcys', '
         passwords'], engine='python')
     ```

     Snipped of passwords containing errors / not loading with pandas (The passwords in this case are whitespaces):

     ```
     4750:      340
     34329:       61
     50383:       44
     99361:       24
     154363:       16
     231236:       11
     562997:        5
     2459758:        2
     2465704:        1
     ```

3. Then we came up with the idea of using preprocessing before reading the file with pandas. There should now be a comma between each password and the password frequency. The runtime has been improved, but there are several 'password errors' again, because many passwords contain commas.

```
1  4750:340,
2  135833:18,,./,./
3  165384:15,,.,.,.
4  210363:12,,,,,,,
5  288009:9,,mnbvcxz
6  380660:7,,,,...
7  453657:6,,fojkiyd
8  453658:6,,f9tovp
9  562846:5,,u8;k,l6-
10 562847:5,,bo9ik
11 562848:5,,b56okpo
12 745987:4,,y]]bdk
13 745988:4,,tgs,ujp;
14 745989:4,,r87hwtwycv
15 745990:4,,kpfNojkiyd
16 745991:4,,kp,bhomN
17 745992:4,,Iuiy9oN
18 745993:4,,./;\\'#
19 745994:4,,.,.,.,.,.
20 844993:3,nan
```

4. In preprocessing we replaced the comma with a tab. Now there are almost no password errors anymore.

```
1  #preprocess txt file and open it with pandas
2  def a():
3      #preprocess the file and save it
4      #only if the file doesent already exists
5      if not os.path.isfile(p_datafile):
6          with open(p_datafile, 'w') as p_file:
7              with open(datafile, 'r') as file:
8                  for line in file:
9                      #add a 'tab' after each passwort frequentcy so it can
   be used as sep
10                     processed_line = re.sub(r'^(\s+)([0-9]+)(\s)',r'\2\t',
   line)
11                     p_file.write(processed_line)

13      df = pd.read_csv(p_datafile, sep='\t', header=None,  names=['frequentcys',
   'passwords'])
14      return df
```

```
1  4750:340
2  844993:3 nan
3  4986233:1 num010975427
4  10807123:1 NAC1993lac
```

(b) *Describe the password set and passwords, e.g. size, distribution, password length, errors in the set. Each description should include at least the center (mean, median or mode), the dispersion (variance, range, percentile).*

- 

```
1  #find, write and correct errors and write some usefull data to file
2  def b(df):
3      #find all NaN/ error rows
4      nan_rows = df[df['passwords'].isna()]
5      #save the index of error rows to list
6      list_nan_index = nan_rows.index.tolist()
7      #write error rows to a new file
8      with open(error_file, 'w') as e_file:
9          for index in list_nan_index:
10             #find line in file
11             line = linecache.getline(p_datafile, index+1)
12             format_line = str(index) + ':' + line
13             e_file.write(format_line)

15      #correct errors in df (drop the rows)
16      df = df.drop(df.index[list_nan_index])
17      #reset index
18      df.reset_index(inplace=True, drop=True)

20      #start writing stats
21      with open(stats_file, 'w') as file:
22          file.write(df.describe().to_string() + '\n')
23          file.write('median:\t' + str(df['frequentcys'].median()) + '\n')
24          file.write('mode:\t' + str(df['frequentcys'].mode()) + '\n' )
25          file.write('var:\t' + str(df.var()) + '\n' )
26      return df
```

```
1           frequentcys
2  count   1.434410e+07
3  mean    2.272893e+00
4  std     9.229819e+01
5  min     1.000000e+00
6  25%     1.000000e+00
7  50%     1.000000e+00
8  75%     1.000000e+00
9  max     2.907290e+05
10 median: 1.0
11 mode: 0    1
12 dtype: int64
13 var: frequentcys    8518.955389
14 dtype: float64
```

(c) *Take a look at the ten most frequent passwords. Give a password (creation) policy for them. Your policy can be informal or very formal.*

- The ten most frequent passwords are the following:

```
1     frequentcys   passwords
2  0       290729      123456
3  1        79076       12345
4  2        76789   123456789
5  3        59462    password
6  4        49952    iloveyou
7  5        33291    princess
8  6        21725     1234567
9  7        20901     rockyou
10 8        20553    12345678
11 9        16648      abc123
```

- These passwords underline the following policy's:

1. 5-9 characters

2. characters can be from the alphabet [0-9] or [a,b-z]

(d) *There is no single password policy. Select all passwords which are 7 to 32 characters long, contain at least one digit and at least one upper case letter. We call this set P1.*

```
1  #create passwordset P1
2  def d_e(df):
3      #\A from beginning of the string till \z end: looks for a word with 7-32 chars
4      a = '\A(?=\w{7,32}\Z)'
5      # [^a-z]* not A-Z 0-many times and then one time A-Z: looks for at least one
       upper case char
6      b = '(?=[^A-Z]*[A-Z])'
7      # D* no digits for 0-many times and then one digit: looks for at least one
       digit
8      c = '(?=\D*\d)'
9      p1 = df.loc[df['passwords'].str.contains(a+b+c, regex=True)].copy()
10     #reset index
11     p1.reset_index(inplace=True, drop=True)
12     return p1
```

(e) *P1 will no contain all passwords which match to our policy. How do we handle passwords which do not appear but match to our policy? You either extend P1 or not. Give the reasons for your decision.*

- We did not extend P1. The strength of a password policy is measured empirically with a large sample of passwords. We want to know how strong the password policy is in a real world scenario. In a real world scenario, the passwords are created by people. Therefore the provided password sample file is ideal, as it contains passwords of accounts from real people. This means that if we were to extend P1 with passwords generated by an algorithm, it would falsify our desired result, "because passwords are not usually distributed uniformly"(quote from Lecture 3)

(f) *Calculate the probability of each password in P1.*

```
1  #calculate the probability of each password in p1
2  def f(p1, sum_of_frequentcys):
3      p1['probability'] = p1['frequentcys'] / sum_of_frequentcys
4      return p1
```

```
1         frequentcys              passwords   probability
2  0              449              Password1      0.000523
3  1              396              PASSWORD1      0.000461
4  2              298              PRINCESS1      0.000347
5  3              236              BABYGIRL1      0.000275
6  4              225              Princess1      0.000262
7  5              217               JORDAN23      0.000253
8  6              207               Passw0rd      0.000241
9  7              203               LOVE123       0.000236
10 8              175               A123456       0.000204
```

(g) *Calculate Shannon Entropy of P1.*

- We calculated the Shannon Entropy with the following formula from the lecture:

$$H(X) = -\sum_{i=1}^{\mathbb{N}} Pr[X = x_i]log_2(Pr[X = x_i])$$

```
1  #calculate Shannon Entropy of p1
2  def g(p1, sum_of_frequentcys):
3      shannon_entropy = (p1['probability'] * numpy.log2(p1['probability'])*-1).sum()
4      return shannon_entropy
```

The result is 19.10637335027745

4

(h) *Assume an online adversary with unlimited guesses likes to break 200 accounts. Estimate the expected number of guesses.*

- For calculating the expected number of guesses we also used a formula from the lecture:

$$G(X) = \sum_{i=1}^{\mathbb{N}} i \cdot Pr[X = x_i]$$

```
1  #calculates expected num of guesses to crack k accounts
2  def h(p1, num_of_accounts_to_crack):
3      #sort values by probability from high to low (i > i+1)
4      p1.sort_values(by=['probability'], ascending = False, inplace = True)
5      #reset index
6      p1.reset_index(inplace=True, drop=True)
7      #calc guesswork
8      guesswork = (p1['probability']* (p1.index +1 )).sum()
9      return num_of_accounts_to_crack*guesswork
```

The result is 57652515

(i) *Assume an online adversary likes to break at least 50 of 200 accounts. Estimate adversary's work.*

- We calculated the Shannon Entropy with the following formula from the lecture:

$$\sum_{i=1}^{\beta} Pr[X = x_i] \geq \frac{l}{k}$$

```
1  #calculates expected num of guesses to crack l accounts of k existing accounts
2  def i(p1):
3      l_div_k = 50/200
4      sum = 0
5      beta = 0
6      while sum < l_div_k:
7          sum += p1['probability'][beta]
8          beta += 1
9      return beta
```

The result is 62693

(j) *Formalize your own password policy and repeat the above estimations.*

- Our passowrd policy's follows the following rules:

  1. 8-64 characters

  2. contains at least two digit

  3. contains at least two upper case letter

  4. contains at least one lower case letter

- Result for the above estimations:

  f) **Probability**

|   | frequentcys | passwords | probability |
|---|---|---|---|
| 0 | 41 | ABCabc123 | 0.000563 |
| 1 | 30 | ABC123abc | 0.000412 |
| 2 | 23 | abc123ABC | 0.000316 |
| 3 | 13 | 123qweASD | 0.000179 |
| 4 | 12 | abcABC123 | 0.000165 |
| 5 | 11 | Sail2Boat3 | 0.000151 |
| 6 | 10 | JamesBond007 | 0.000137 |
| 7 | 10 | 2Bornot2B | 0.000137 |
| 8 | 9 | ABCdef123 | 0.000124 |

  g) **Shanon Entropy** result is 16.036519282236434

  h) The adversary needs 6565494.616061121 guesses to **break 200 accounts**

  i) The adversary needs 14473 guesses to **break 50 of 200 accounts**

(k) *Compare the password policies and estimations.*

The main difference between p1 and our own password policy is that we require two capital letters and and two numbers. We thought that this would make the passwords created by users more secure. However, the Shannon entropy and the number of guesses for cracking accounts have decreased. Perhaps this is due to the lower password sample size, or it is because our password policy restricts the number of possible passwords and thus becomes weaker. However, we tend to assume the first.

| Estimations | P1 | Own password policy |
|---|---|---|
| Password sample size | 858.415 | 72.772 |
| Shanon Entropy | 19 | 16 |
| Guesses to break 200 accounts | 57.652.515 | 6.565.494 |
| Guesses to break 50 of 200 accounts | 62.693 | 14.473 |

# 2 Project

(a) *Setup Docker on your machine (or on a virtual machine).*

(b) *Add a docker image for your project. You find an example in the resources folder.*

a) *Some use full docker commands:*

```
1  #lists all images
2  docker image ls
3  #lists all containers
4  docker ps -a

6  #stops all containers
7  docker stop $(docker ps -a -q)
8  #deletes all containers
9  docker rm $(docker ps -a -q)
10 #deletes all images
11 docker rmi --all --force

13 #builds images of docker_compose.yml
14 docker-compose build
15 #creates containers of docker_compose.yml in detached mode and overwrites old
       images
16 docker-compose up -d --force-recreate
17 #deletes all images and containers created by last docker-compose up call
18 docker-compose down
```

b) *Docker and docker-compose files:*

```
1  FROM python:3.8-alpine
2  ADD client.py /client/
3  WORKDIR /client/
```

```
1  FROM python:3.8-alpine
2  ADD client.py /client/
3  WORKDIR /client/
```

```
1   # file from https://www.freecodecamp.org/news/a-beginners-guide-to-docker-how-to-create-a-
        client-server-side-with-docker-compose-12c8cf0ae0aa/
2   # A docker-compose must always start by the version tag.
3   # We use "3" because it's the last version at this time.
4   version: "3"

6   # You should know that docker-composes works with services.
7   # 1 service = 1 container.
8   # For example, a service maybe, a server, a client, a database...
9   # We use the keyword 'services' to start to create services.
10  services:
11    # First service (container): the server.
12    # Here you are free to choose the keyword.
13    # It will allow you to define what the service corresponds to.
14    # We use the keyword 'server' for the server.
15    server:
16      # The keyword "build" will allow you to define
17      # the path to the Dockerfile to use to create the image
18      # that will allow you to execute the service.
19      # Here 'server/' corresponds to the path to the server folder
20      # that contains the Dockerfile to use.
21      build: server/

23      # The command to execute once the image is created.
24      # The following command will execute "python ./server.py".
25      command: python ./server.py2

27      # Remember that we defined in'server/server.py' 1234 as port.
28      # If we want to access the server from our computer (outside the container),
29      # we must share the content port with our computer's port.
30      # To do this, the keyword 'ports' will help us.
31      # Its syntax is as follows: [port we want on our machine]:[port we want to retrieve in the
          container]
32      # In our case, we want to use port 1234 on our machine and
33      # retrieve port 1234 from the container (because it is on this port that
34      # we broadcast the server).
35      ports:
36        - 1234:1234

38    # Second service (container): the client.
39    # We use the keyword 'client' for the server.
40    client:
41      # docker run -i
42      # activates std input
43      stdin_open: true
44      # docker run -t
45      # activates std output
46      tty: true

48      # Here 'client/ corresponds to the path to the client folder
49      # that contains the Dockerfile to use.
50      build: client/

52      # The command to execute once the image is created.
53      # The following command will execute "python ./client.py".
54      command: python ./client.py

56      # The keyword 'network_mode' is used to define the network type.
57      # Here we define that the container can access to the 'localhost' of the computer.
58      network_mode: host

60      # The keyword'depends_on' allows you to define whether the service
61      # should wait until other services are ready before launching.
62      # Here, we want the 'client' service to wait until the 'server' service is ready.
63      depends_on:
64        - server

66      # starts this service as container 'scale' times
67      scale: 3
```

(c) Before using third party dependencies. We might should talk about them and potential threats. Read
(or watch the video) the paper *Small World with High Risks: A Study of Security Threats in the npm
Ecosystem* and list at least 3 risks using npm. TODO

(d) Is this relevant for Python and our project? Give a short statement. You might find some research
about Python, too. TODO